



A privacy self-assessment framework for online social networks



Ruggero G. Pensa*, Gianpiero Di Blasi

Dept. of Computer Science, University of Torino, C.So Svizzera, 185 – I-10149 Torino, Italy

ARTICLE INFO

Article history:

Received 19 December 2016
Revised 11 April 2017
Accepted 20 May 2017
Available online 22 May 2017

Keywords:

Privacy measures
Online social networks
Active learning

ABSTRACT

During our digital social life, we share terabytes of information that can potentially reveal private facts and personality traits to unexpected strangers. Despite the research efforts aiming at providing efficient solutions for the anonymization of huge databases (including networked data), in online social networks the most powerful privacy protection “weapons” are the users themselves. However, most users are not aware of the risks derived by the indiscriminate disclosure of their personal data. Moreover, even when social networking platforms allow their participants to control the privacy level of every published item, adopting a correct privacy policy is often an annoying and frustrating task and many users prefer to adopt simple but extreme strategies such as “visible-to-all” (exposing themselves to the highest risk), or “hidden-to-all” (wasting the positive social and economic potential of social networking websites). In this paper we propose a theoretical framework to i) measure the privacy risk of the users and alert them whenever their privacy is compromised and ii) help the users customize semi-automatically their privacy settings by limiting the number of manual operations. By investigating the relationship between the privacy measure and privacy preferences of real Facebook users, we show the effectiveness of our framework.

© 2017 Elsevier Ltd. All rights reserved.

1. Introduction

Social networks are one of the main traffic sources in the Internet. At the end of 2014, they attracted more than 31% of the worldwide Internet traffic towards the Web. Facebook, the most famous social networking platform, drives alone 25% of the whole traffic. As a comparison, Google search engine represents just over 37% of the global traffic. More than two billions people are estimated to be registered in at least one of the most popular social media platforms (Facebook hits the goal of one billion users in 2012). Overall, the number of active “social” accounts are more than two billions. In view of these numbers, the risks due to a more and more global and unaware diffusion of our sensitive and less sensitive personal data cannot be overlooked. If, on the one hand, many users are informed about the risks linked to the disclosure of personal facts (private life events, sexual preferences, diseases, political ideas, and so on), on the other hand the awareness of being exposed to privacy breaches each time we disclose facts that are apparently less sensitive is still insufficiently widespread. A GPS tag far from home or pictures taken during a journey, may alert potential thieves who may clean out the apartment. The disclosure of family relation-

ships may expose our own or other family members' privacy to the risks of stalking, slander and cyberbullying. Moreover, the research project myPersonality (Kosinski, Stillwell, & Graepel, 2013), carried out at the University of Cambridge, has shown that, by leveraging Facebook user's activity (such as “Likes” to posts or fan pages) it is possible to “guess” some very private traits of the user's personality. According to another study, it is even possible to infer some user characteristics from the attributes of users who are part of the same communities (Mislove, Viswanath, Gummadi, & Druschel, 2010). As a consequence, privacy has become a primary concern among social network analysts and Web/data scientists. Also, in recent years, many companies are realizing the necessity to consider privacy at every stage of their business. In practice, they have been turning to the principle of *Privacy by Design* (Cavoukian, 2012) by integrating privacy requirements into their business model.

Despite the huge research efforts aiming at providing efficient solutions to the anonymization of huge databases (including networked data) (Backstrom, Dwork, & Kleinberg, 2011; Xue, Karras, Raïssi, Kalnis, & Pung, 2012; Zhou & Pei, 2011; Zou, Chen, & Özsu, 2009), in online social networks the most powerful privacy protection is in the hands of the users: they, and only they, decide what to publish and to whom. Even though social networking sites (such as Facebook), notify their users about the risks of disclosing private information, most people are not aware of the dangers due to the indiscriminate disclosure of their personal data. Moreover, despite the fact that all social media provide some advanced tools for con-

* Corresponding author.

E-mail addresses: ruggero.pensa@unito.it (R.G. Pensa), diblasi@di.unito.it (G. Di Blasi).

trolling the privacy settings of the user's profile, such tools are not user-friendly and they are barely utilized, in practice. According to Facebook former CTO Bret Taylor, most people have modified their privacy settings, but in 2012, still "13 million users [in the United States] said they had never set, or didn't know about, Facebook's privacy tools". Often the choices of many users are limited to two: i) make their own profile completely public, being exposed to all the above mentioned risks, ii) make their own profile completely private, preventing all opportunities offered by the social network sites. Some studies try to foster risk perception and awareness by "measuring" users' profile privacy according to their privacy settings (Liu & Terzi, 2010; Wang, Nepali, & Nikolai, 2014). These metrics usually require a *separation-based* policy configuration: in other terms, the users decide "how distant" a published item may spread in the network. Typical separation-based privacy policies for profile item/post visibility include: visible to no one, visible to friends, visible to friends of friends, public. However, this policy fails when the number of user friends becomes large. According to a well-known anthropological theory, in fact, the maximum number of people with whom one can maintain stable social (and cybersocial) relationships (known as Dunbar's number) is around 150 (Dunbar, 2016; Roberts, Dunbar, Pollet, & Kuppens, 2009), but the average number of user friends in Facebook is more than double¹. This means that many social links are weak (offline and online interactions with them are sporadic), and a user who sets the privacy level of an item to "visible to friends" probably is not willing to make that item visible to *all* her friends. Other studies try to make the customization process of the privacy settings less frustrating (Fang & LeFevre, 2010). However, a consensus on how to identify a trade-off between privacy protection and exploitation of social network potentials is still far from being achieved.

With the final goal of enhancing users' privacy awareness in online social networks, in this paper we propose a theoretical framework to i) measure the privacy risk of the users and alert them whenever their privacy is compromised and ii) help the exposed users customize semi-automatically their privacy level by limiting the number of manual operations thanks to an active learning approach. Moreover, instead of using a *separation-based* policy for computing the privacy risk, in this paper we adopt a *circle-based* formulation of the privacy score proposed by Liu and Terzi (2010). We assume that a user may set the visibility of each action and profile item separately for each other user in her friend list. For instance, a user u may decide to allow the access to all photo albums to friends f_1 and f_2 , but not to friend f_3 . In our score, the sensitivity and visibility of profile item i published by user u are computed according to the set of u 's friends that are allowed to access the information provided by i . We show experimentally that our circle-based definition of privacy score better capture the real privacy leakage risk. Moreover, by investigating the relationship between the privacy measure and the privacy preferences of real Facebook users, we show that our framework may effectively support a safer and more fruitful experience in social networking sites. Differently from other research works addressing the same problem, our framework takes into account both users' preferences and the real sensitive information leakage risk in deciding how much visibility should be given to each profile item.

Our contribution can be resumed as follows:

- we define a formal framework for privacy self-assessment in online social networks based on both sensitivity and visibility of user profile items;
- we use a new privacy score leveraging more accurate *circle-based* policies;

- we present a semi-supervised machine learning approach to support the configuration of the visibility level of user profile items;
- we report the results of several experiments on original data obtained from real Facebook users.

The remainder of the paper is organized as follows: we briefly review the related literature in Section 2; the overview and the theoretical details of our framework are presented in Section 3; Section 4 provides the report of our experimental validation; finally, we draw some conclusions, discuss some limitations and propose some future research directions in Section 5.

2. Related work

With the unrestrained success of online social networks, there has been increasing research interests about privacy protection methods for individuals that participate in them. Most research efforts are devoted to the identification and formalization of privacy breaches and to the anonymization of networked data. The goal is to modify data so that the probability of identifying an individual within the network is minimized. This objective is achieved by either anonymizing only the network structure or anonymizing both network structure and user attributes (Zheleva & Getoor, 2011).

Some of the most relevant contributions tackle the problem of graph anonymization by applying edge modification (Liu & Terzi, 2008; Zhou & Pei, 2011; Zou, Chen, & Özsu, 2009), randomization (Vuokko & Terzi, 2010; Ying & Wu, 2011), generalization (Cormode, Srivastava, Bhagat, & Krishnamurthy, 2009; Hay, Miklau, Jensen, Towsley, & Weis, 2008) or differentially private mechanisms (Hay, Li, Miklau, & Jensen, 2009; Task & Clifton, 2012). Among the approaches that anonymize also the user attributes, Zhou and Pei (2011) adopt a greedy edge modification and label generalization algorithm, Zheleva and Getoor (2008) anonymize nodes attribute first and then tries to preserve the network structure, Campan and Truta (2009) optimize an utility function using the attribute and structural information simultaneously.

All these works focus on how to share social networks owned by companies or organizations masking the identities or the sensitive connections of the individuals involved. However, less attention has been given to the privacy risk of users caused by their information-sharing activities (e.g., posts, likes, shares). In fact, since disclosing information on the web is a voluntary activity, a common opinion is that users should care about their privacy and control it during their interaction with other social network users. Although multiple complex factors are involved in user privacy protection on social media (Litt, 2013), privacy controls for online social networking sites are not fully socially aware (Misra & Such, 2016) and are barely utilized in practice. This statement is confirmed by a study of Liu, Gummadi, Krishnamurthy, and Mislove (2011) which shows that 36% of Facebook content is shared with the default privacy settings and exposed to more users than expected.

Thus, another branch of research has focused on investigating measures, strategies and tools to enhance the users' privacy awareness and help them act more safely during their day-to-day social network activity. Liu and Terzi (2010) propose a framework to compute a privacy score measuring the users' potential risk caused by their participation in the network. This score takes into account the sensitivity and the visibility of the disclosed information and leverages the item response theory as theoretical basis for the mathematical formulation of the score. Instead, Motahari, Zivarras, and Jones (2010) propose an information-theoretic estimation of the user anonymity level to help predict the identity inference risks according to both external knowledge and the correlation between user attributes. Cetto et al. (2014) present an online

¹ <http://www.pewresearch.org/fact-tank/2014/02/03/6-new-facts-about-facebook/>

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات